

Gödel for Goldilocks: A Rigorous, Streamlined Proof of (a variant of) Gödel’s First Incompleteness Theorem¹

Dan Gusfield

Department of Computer Science, UC Davis

August 2014, revised November 15, 2014

1 Introduction: Why I wrote this

Gödel’s famous incompleteness theorems (there are two of them) concern the ability of a formal system to state and derive all true statements, and only true statements, in some fixed domain; and concern the ability of logic to determine if a formal system has that property. They were developed in the early 1930s. Very loosely, the first theorem says that in any “sufficiently rich” formal proof system where it is not possible to prove a false statement about *arithmetic*, there will also be true statements about arithmetic that cannot be proved.

Most discussions of Gödel’s theorems fall into one of two types: either they emphasize perceived cultural and philosophical meanings of the theorems, and perhaps sketch some of the ideas of the proofs, usually relating Gödel’s proofs to riddles and paradoxes, but do not attempt rigorous, complete proofs; or they do present rigorous proofs, but in the traditional style of mathematical logic, with all of its heavy notation, difficult definitions, technical issues in Gödel’s original approach, and connections to broader logical theory before and after Gödel. Many people are frustrated by these two extreme types of expositions² and want a short, straight-forward, rigorous proof that they can understand.

Over time, various people have realized that somewhat weaker, but still meaningful, variants of Gödel’s first incompleteness theorem can be rigorously proved by simpler arguments based on notions of computability. This approach avoids the heavy machinery of mathematical logic at one extreme; and does not rely on analogies, paradoxes, philosophical discussions or hand-waiving, at the other extreme. This is the just-right *Goldilocks* approach. However, the available expositions of this middle approach have

¹This exposition requires minimal background. Other than common things, you only need to know what an integer is; what a function is; and what a computer program is. You do need several hours, and you need to focus—the material is concrete and understandable, but it is not trivial. This material was the basis for the first two lectures of my course offering “The Theory of Computation” (a sophomore/junior level course) in October 2014.

²To verify this, just randomly search the web for questions about Gödel’s theorem.

still been aimed at a relatively sophisticated audience, and have either been too brief,³ or have been embedded in larger, more involved discussions.⁴ A short, self-contained Goldilocks exposition of a version of Gödel's first theorem, aimed at a broad audience, has been lacking. Here I offer such an exposition.

2 There are Non-Computable Functions

We start with a discussion of computable and non-computable functions.

Definition offer We use Q to denote all functions from the positive integers to $\{0, 1\}$. That is, if f is in Q , then for any positive integer x , $f(x)$ is either 0 or 1.

Note that since a function in Q is defined on *all* positive integers, the number of functions in Q is infinite.

Definition Define a function f in Q to be *computable* if there is a finite-sized computer program (in Python, for example) that executes on a computer (a MacBook Pro running Snow Leopard, for example) that computes function f . That is, given *any* positive integer x , the program finishes in finite time and correctly spits out the value $f(x)$.

Definition Let A be the set of functions in Q that are computable.

Note that the number of functions in A is infinite. For example, the function $f(7) = 1$ and $f(x) = 0$ for all $x \neq 7$ is a computable function, and we can create a similar computable function for any positive integer, in place of 7. So, since there are an infinite number of positive integers, there are an infinite number of computable functions.

Theorem 2.1 *There are functions in Q that are not computable. That is, $A \subset Q$.*

Proof In this proof we would like to talk about an *ordering* (or an ordered list)⁵ of all functions in A , rather than just the *set* A . It might seem self-evident that such an ordering should exist, and so one might think we could just assert that it does. But issues of ordering are subtle; there are unsettled questions about which properties are sufficient to guarantee that an ordering exists. So, we want to be careful and fully establish that an ordering of the functions in A does exist.⁶

³For example, in Scott Aaronson's book *Quantum Computing Since Democritus*.

⁴For example, Sipser's excellent book on the Theory of Computation, where the exposition of Gödel's theorem relies on an understanding of Turing machines and the Undecidability of the Halting problem. Another example is *An Introduction to Gödel's Theorems* by Peter Smith, which develops much more logical machinery before proving a variant of Gödel's theorem. But, for anyone wanting a readable, deeper and broader treatment of the theorems than I present here, I highly recommend that book.

⁵The common notion of an ordering is more technically called a *total order*.

⁶I thank the students in CS 120 Fall 2014 who asked why an earlier draft of this exposition goes into such detail on the existence of an ordering.

An ordering Exists: First, choose a computer language and consider a program in that language. Each line in a program has some end-of-line symbol, so we can concatenate the lines together into a single long string. Therefore, we think of a program in that computer language as a single string written using some finite alphabet.

Now, since A consists of the computable functions in Q , for each function $f \in A$, there is some computer program P_f (in the chosen computer language) that computes f . Program P_f (considered as a single string) has some finite length. We can, conceptually, order the strings representing the programs that compute the functions in A into a list L in *order* of the lengths of the strings. To make the ordering perfectly precise, when there are strings of the same length, we order those strings lexicographically (i.e., the way they would be *alphabetically* ordered in a dictionary). So, each *program* that computes a function in A has a *well-defined* position in L . Then, since each function in A is computed by some program in the ordered list L , L also defines an ordered list, which we call L' , containing all the functions in A .

A function f in A might be computed by different computer programs, so f might appear in L' more than once. If that occurs, we could, conceptually, remove all but the *first* occurrence of f in L' , resulting in an ordering of the functions in A , as desired. We will see that it will not harm anything if f is computed by more than one program in L , and hence appears in L' more than once. The only point that will matter is that there is some ordered listing L' of the functions in A that includes every function in A .

Let f_i denote the function in A that appears in position i in L' ; that is, f_i is computed by the i 'th program in L . (Remember that lists L and T are only conceptual; we don't actually build them—we only have to imagine them for the sake of the proof). Next, consider a table T with one column for each positive integer, and one row for each program in L ; and associate the function f_i with row i of T . Then set the value of cell $T(i, x)$ to $f_i(x)$. See Table 1.

Function \bar{f} : Next, we define the function \bar{f} from the positive integers to $\{0, 1\}$ as $\bar{f}(i) = 1 - f_i(i)$. For example, based on the functions in Table 1, $\bar{f}(1) = 0$; $\bar{f}(2) = 1$; $\bar{f}(3) = 1$; $\bar{f}(4) = 0$; $\bar{f}(5) = 1$

Note that in the definition of $\bar{f}(i)$, the same integer i is used both to identify the function f_i in A , and as the input value to f_i and to \bar{f} . Hence the values for \bar{f} are determined from the values along the main *diagonal* of table T . Note also that \bar{f} changes 0 to 1, and changes 1 to 0. So, the values of function \bar{f} are the *opposite* of the values along the main diagonal of Table T . Clearly, function \bar{f} is in Q .

Now we ask: Is \bar{f} a computable function?

The answer is no for the following reason. If \bar{f} were a computable function, then there would be some row i^* in T such that $\bar{f}(x) = f_{i^*}(x)$ for every positive integer x . For example, maybe i^* is 57. But $\bar{f}(57) = 1 - f_{57}(57) \neq f_{57}(57)$, so \bar{f} can't be f_{57} . More generally, $\bar{f}(i^*) = 1 - f_{i^*}(i^*)$, so \bar{f} and f_{i^*} differ at least for one input value (namely i^*),

	1	2	3	4	5	.	x	.	i	...
f_1	1	1	0	0	0		$f_1(x)$			
f_2	0	0	1	0	0		$f_2(x)$			
f_3	1	1	0	0	1		$f_3(x)$			
f_4	0	0	1	1	0		$f_4(x)$			
f_5	0	1	0	0	0		$f_5(x)$			
.						.				
.							.			
.								.		
f_i									$f_i(i)$	
\vdots										\ddots

Table 1: The conceptual Table T is contains an ordered list of all computable functions in Q , and their values at all of the positive integers.

so $\bar{f} \neq f_{i^*}$. Hence, there is no row in T corresponding to \bar{f} , and so \bar{f} is not in set A . So \bar{f} is *not* computable—it is in Q , but not in A . ■

3 What is a Formal Proof System?

How do we connect Theorem 2.1, which is about functions, to Gödel’s first incompleteness theorem, which is about logical systems? We first must define a *formal proof system*.

Definition A formal proof system Π has three components:

1. A finite alphabet, and some finite subset words and phrases that can be used in forming (or writing) *statements*.⁷
2. A finite list of *axioms* (statements that we take as true); and
3. A finite list of *rules of reasoning*, also called *logical inference*, *deduction* or *derivation* rules, that can be applied to create a new statement from axioms and the statements already created, in an unambiguous, mechanical way.

The word “mechanical” is central to the definition of rules of reasoning, and to the whole purpose of a formal proof system:

⁷These words and phrases are strings in the alphabet of Π . We will say that they are a subset of English.

... we need to impose some condition to the effect that recognizing an axiom or applying a rule must be a mechanical matter ... it is required of a formal system that in order to verify that something is an axiom or an application of a rule of reasoning, we ... need only apply mechanical checking of the kind that can be carried out by a computer.⁸

For example, the alphabet might be the standard ASCII alphabet with 256 symbols, and Axiom 1 might be: “for any integer x , $x + 1 > x$.” Axiom 2 might be: “for any integers x and y , $x + y$ is an integer.” A derivation rule might be: “for any three integers, x, y, z , if $x > y$ and $y > z$ then $x > z$.” (Call this rule the “Transitivity Rule”.)

The finite set of allowed English words and phrases might include the phrase: “for any integer”. Of course, there will typically be more axioms, derivation rules, and known words and phrases than in this example.

3.1 What is a Formal Derivation?

Definition A *formal derivation* in Π of a statement S is a series of statements that begin with some axioms of Π , and then successively apply derivation rules in Π to obtain statement S .

For example, S might be the statement: “For any integer w , $w + 1 + 1 > w$ ”. A formal derivation of S in Π (using axioms and derivation rules introduced above) might be:

- w is an integer, 1 is an integer, so $w + 1$ is an integer (by Axiom 2).
- $w + 1$ is an integer (by the previous statement), 1 is an integer, so $w + 1 + 1 > w + 1$ (by Axiom 1).
- $w + 1 + 1 + 1 > w + 1 > w$ (by the previous statement and Axiom 1).
- $w + 1 + 1 > w$ (by the Transitivity Rule). This is statement S .

The finite subset of English used in this formal derivation includes the words and phrases “is an integer”, “by the Transitivity Rule”, “by the previous statement” etc. . These would be part of the finite subset of English that is part of the definition of Π . Each phrase used must have a clear and precise meaning in Π , so that each statement in a formal derivation, other than an axiom, follows in a mechanical way from the preceding statements by the application of some derivation rule(s) or axioms.

Formal derivations are very tedious, and humans don’t want to write derivations this way, but computers can write and check them, a fact that is key in our treatment of

⁸Gödel’s Theorem, by Torkel Franzen, CRC Press, 2005.

Gödel's theorem. (Note that what I have called a “formal derivation” is more often called a “formal proof”. But that is confusing, because people usually think of a “proof” as something that establishes a *true* statement, not a statement that might be false. So here we use “formal derivation” to avoid that confusion.)

3.2 Mechanical Generation and Checking of Formal Derivations

We now make four key points about formal derivations.

1. It is easy to write a program P that can begin generating, in order of the lengths of the strings, every string s that can be written in the alphabet of Π , and using allowed words and phrases of the formal proof system Π . Program P will never stop because there is no bound on the length of the strings, and most of the strings will not be formal derivations of anything. But, for any finite-length string s using the alphabet of Π , P will eventually (and in finite time) generate s .

2. A formal derivation, being a series of statements, is just a string formed from the alphabet and the allowed words and phrases of the formal proof system Π . Hence, if s is any string specifying a formal derivation, P will eventually (and in finite time) generate it.

3. We can create a program P' that knows the alphabet, the axioms, the deduction rules, and the meaning of the words of the allowed subset of English used in Π , so that P' can precisely interpret the effect of each line of a formal derivation. That is, P' can *mechanically* check whether each line is an axiom, or follows from the previous lines by an application of some deduction rule(s) or axioms. Therefore, given a statement S , and a string s that might be a formal derivation of S , program P' can check (in a purely mechanical way, and in finite time) whether string s is a formal derivation of statement S in Π .

4. For any statement S , after program P generates a string s , program P' can check whether s is a formal derivation of statement S in Π , before P generates the next string. Hence, if there is a formal derivation s in Π of statement S , then s will be generated and recognized in finite time by interleaving the execution of programs P and P' .

Note that most of the strings that P generates will be garbage, and most of the strings that are not garbage will not be formal derivations of S in Π . But, if string s is a formal derivation of statement S , then in finite time, program P will generate s , and program P' will recognize that s is a formal derivation in Π of statement S .

Similarly, we can have another program P'' that checks whether a string s is a formal derivation of the statement “not S ”, written $\neg S$. So if $\neg S$ is a statement that can be derived in Π , the interleaved execution of programs P and P'' will, in finite time, generate and recognize that s is a formal derivation of $\neg S$.

4 Back to Gödel

How do we connect all this to Gödel's first incompleteness theorem? We want to show the variant of Gödel's theorem that says: in any “rich-enough” formal proof system where no false statement about functions can be derived, there are true statements about functions that cannot be derived. We haven't defined what “true” or “rich-enough” means in general, but we will in a specific context.

Recall function \bar{f} , and recall that it is well-defined, i.e., there is a value $\bar{f}(x)$ for every positive integer x , and for any specific x , $\bar{f}(x)$ is either 0 or 1. Recall also, that \bar{f} is not a computable function.

Definition We call a statement an \bar{f} -statement if it is either:

“ $\bar{f}(x)$ is 1,”

or:

“ $\bar{f}(x)$ is 0,”

for some positive integer x .

Note that every \bar{f} -statement is a statement about a specific integer. For example the statement “ $\bar{f}(57)$ is 1” is an \bar{f} -statement, where x has the value 57. Since, for any positive integer x , $\bar{f}(x)$ has only two possible values, 0 or 1, when the two kinds of \bar{f} -statements refer to the same x , we refer to the first statement as $Sf(x)$ and the second statement as $\neg Sf(x)$.

What is Truth? We say an \bar{f} -statement $Sf(x)$ is “true”, and $\neg Sf(x)$ is “false”, if in fact $\bar{f}(x)$ is 1. Similarly, we say an \bar{f} -statement $\neg Sf(x)$ is true, and $Sf(x)$ is false, if in fact $\bar{f}(x)$ is 0. Clearly, for any positive integer x , one of the statements $\{Sf(x), \neg Sf(x)\}$ is true and the other is false. In this context, truth and falsity are simple concepts (not so simple in general).

Clearly, it is a desirable property of a formal proof system Π , that it is not possible to give a formal derivation in Π for a statement that is false.

What does it mean to be rich-enough? We need a definition.

Definition We define a formal proof system Π to be *rich-enough* if any \bar{f} -statement can be formed (i.e., stated, or written) in Π .

Note that the words “formed”, “stated”, “written” do not mean “derived”. The question of whether a statement can be derived in Π is at the heart of Gödel's theorem. Here, we are only saying that the statement can be formed (or written) in Π .

4.1 The Proof of our variant of Gödel's Theorem

Now let Π be a rich-enough formal proof system, and suppose **a**) that Π has the properties that no false \bar{f} -statements can be derived in Π ; and suppose **b**) that for any true \bar{f} -statement S , there is a formal derivation s of S in Π .

Since Π is rich-enough, for any positive integer x , both statements $Sf(x)$ and $\neg Sf(x)$ can be formed in Π , and since exactly one of those statements is true, suppositions **a** and **b** imply that there is a formal derivation in Π of exactly one of the two statements, in particular, the statement that is true. But this leads to a contradiction of the established fact that function \bar{f} is not computable.

In more detail, if the two suppositions (**a** and **b** hold, the following approach describes a computer program P^* that can correctly determine the value of $\bar{f}(x)$, for any positive integer x , in finite time.

Program P^* : Given x , start program P to successively generate all possible strings (using the finite alphabet and known words and phrases in Π), in order of their lengths, breaking ties in length lexicographically (as we did when discussing list L). After P generates a string s , run program P' to see if s is a formal derivation of statement $Sf(x)$. If it is, output that $\bar{f}(x) = 1$ and halt; and if it isn't, run P'' to see if s is a formal derivation of $\neg Sf(x)$. If it is, output that $\bar{f}(x) = 0$ and halt; and if it isn't, let P go on to generate the next possible string.

The two suppositions **a** and **b** guarantee that for any positive integer x , this mechanical computer program, P^* , will halt in finite time, outputting the correct value of $\bar{f}(x)$. But then, \bar{f} would be a *computable* function (computable by program P^*), contradicting the already established fact that \bar{f} is not a computable function. So, the two suppositions **a** and **b** lead to a contradiction, so they cannot both hold for any rich-enough formal proof system Π . There are several equivalent conclusions that result. One is:

Theorem 4.1 *For any rich-enough formal proof system Π in which no formal derivation of a false \bar{f} -statement is possible, there will be some true \bar{f} -statement that cannot be formally derived in Π .*

A different, but equivalent conclusion is:

Theorem 4.2 *In any rich-enough formal proof system Π in which no formal derivation of a false \bar{f} -statement is possible, there will be some positive integer x such that neither statement $Sf(x)$ nor statement $\neg Sf(x)$ can be formally derived.*

We leave to the reader the proof that Theorems 4.2 and 4.1 are equivalent. Theorems 4.1 and 4.2 are variants of Gödel's first incompleteness theorem.

5 A Little Terminology of Formal Logic

We proved Theorem 4.1 with as little terminology from formal logic as possible. That was one of the goals of this exposition. Still, it is useful to introduce some terminology to be more consistent with standard use.

Definition A formal proof system Π is called *sound* if only true statements can be derived in Π . But note that we don't require that *all* true statements be derived in Π .

Definition A formal proof system Π is called *complete* if for any statement S that can be formed in Π , one of the statements S or $\neg S$ can be derived in Π . But note that we don't require that the derived statement be true.

Theorem 4.2 can then be stated as:

Theorem 5.1 *No formal proof system Π that can form any \overline{f} statement can be both sound and complete.*

6 Gödel's Second Incompleteness Theorem

Definition A formal proof system Π is called *consistent* if it is never possible to derive a statement S in Π and also derive the statement $\neg S$ in Π .

Later in the course, we will talk about Gödel's second incompleteness theorem, which needs more machinery. Informally, it says that if Π is rich-enough and consistent, there cannot be a formal derivation in Π of the statement: " Π is consistent". More philosophically, but not precisely, for any (rich enough) formal proof system Π that is consistent, the consistency of Π can only be established by a different formal proof system Π' . (But then, what establishes the consistency of Π' ?)

7 Optional Homework Questions:

1. In two places in the proofs, ties in the lengths of strings are broken lexicographically. I claim that this detail is not needed in either place. Is this true?
2. In the proof of Theorem 2.1, what is the point of requiring the computer programs to be listed in order of their lengths? Would the given proof of Theorem 2.1 remain correct if the programs were (somehow) listed in no predictable order?
3. In program P^* , what is the point of requiring program P to generate strings in order of their lengths? Would the given proof of Theorem 4.1 remain correct if P did not generate the strings in that order, but could (somehow) generate all the strings in no predictable order?

4. Doesn't the following approach show that $\bar{f}(x)$ is computable?

First, create a computer program P' that can look at a string s over the finite alphabet used for computer programs (in some fixed computer language, for example, C), and determine if s is a legal computer program that computes a function f in Q . Certainly, a compiler for C can check if s is a syntactically correct program in C.

Then given any positive integer x , use program P to generate the strings over the finite alphabet used for computer programs, in order of their length, and in the same order as used in table T . After each string s is generated, use program P' to determine if s is a program that computes a function in Q . Continue doing this until x such programs have been found. In terms of table T , that program, call it F , will compute function f_x . Program F has finite length, so P will only generate a finite number of strings before F is generated. Then once F is generated, run it with input x . By definition, program F will compute $f_x(x)$ in finite time. Then output $\bar{f}(x) = 1 - f_x(x)$.

So, this approach seems to be able to compute $\bar{f}(x)$ in finite time, for any positive integer x , showing that \bar{f} is a computable function. Doesn't it?

Discuss and resolve.

5. Use the resolution to the issue in problem 4, to state and prove an interesting theorem about computer programs (yes, this is a vague question, but the kind that real researchers face daily).

6. Show that theorems 4.1 and 4.2 are equivalent.

7. Show that a formal proof-system that is sound is also consistent. Then ponder whether it is true that any formal proof-system that is consistent must be sound. Hint: no.

8 What is not in this exposition?

Lots of stuff that you might see in other proofs and expositions of Gödel's first incompleteness theorem: propositional and predicate logic, models, WFFs, prime numbers, prime factorization theorem, Chinese remainders, Gödel numbering, countable and uncountable infinity, self-reference, recursion, paradoxes, liars, barbers, librarians, "This statement is false", Peano postulates, Zermelo-Fraenkel set theory, Hilbert, Russell, Turing, universal Turing machines, the halting problem, undecidability, ..., quantum theory, insanity, neuroscience, the mind, zen, self-consciousness, evolution, relativity, philosophy, religion, God, Stalin, The first group of topics are actually related in precise, technical ways to the theorem, but can be avoided, as done in this exposition.⁹ Some of those related technical topics are important in their own right, particularly Turing

⁹Most expositions of Gödel's theorems use self-reference, which I find unnecessarily head-spinning, and I think its use is sometimes intended to make Gödel's theorem seem deeper and more mystical than it already is.

undecidability, which we will cover in detail later in the course. The second group of topics are not related in a precise, technical way to the theorem. Some are fascinating in their own right, but their inclusion makes Gödel's theorem more mystical, and should not be confused for its actual technical content.

9 Final Comments

The exposition here does not follow Gödel's original proof, and while the exposition is my own, the general approach reflects (more and less) the contemporary computer-sciencey way that Gödel's theorem is thought about, i.e., via computability. In coming to this exposition for undergraduates, I must acknowledge the discussion of Gödel's theorems in Scott Aaronson's book *Quantum Computing Since Democritus*, and an exposition shown to me by Christos Papadimitriou. Those are both shorter, aimed at a more advanced audience, and are based on the undecidability of Turing's Halting problem. I also thank David Doty for pointing out an incorrect definition in my first version of this exposition.

Second, I must state again that the variant of Gödel's theorem proved here is weaker than what Gödel originally proved. In this exposition, the formal proof system must be able to express any \bar{f} -statement, but Gödel's original proof only requires that the formal system be able to express statements about arithmetic (in fact, statements about arithmetic on integers only using addition and multiplication). This is a more limited domain, implying a stronger theorem. That difference partly explains why a proof of Gödel's original theorem is technically more demanding than the exposition here. Further, Gödel did not just prove the *existence* of a true statement that could not be derived (in any sufficiently rich, sound proof system), he demonstrated a *particular* statement with that property. But, I believe that the moral, cultural, mathematical, and philosophical impact of the variant of Gödel's theorem proved here is comparable to that of Gödel's actual first incompleteness theorem. Many modern treatments of Gödel's theorem similarly reflect this view. Of course, some people disagree and insist that anything using the phrase "Gödel's theorem" must actually be the same as what Gödel proved.¹⁰

¹⁰Although, Gödel did not actually prove what is generally stated as "Gödel's theorem", but only a weaker form of it—which was later strengthened by Rosser to become the classic "Gödel's theorem". Accordingly, some people call it the "Gödel-Rosser theorem".